



# On $C^0$ -persistent homology and trees

Daniel Perez

## ► To cite this version:

| Daniel Perez. On  $C^0$ -persistent homology and trees. 2020. hal-03040819v2

**HAL Id: hal-03040819**

**<https://hal.science/hal-03040819v2>**

Preprint submitted on 7 Dec 2020 (v2), last revised 23 May 2022 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# On $C^0$ -persistent homology and trees

Daniel Perez<sup>\*1,2,3</sup>

<sup>1</sup>Département de mathématiques et applications, École normale supérieure, CNRS, PSL University, 75005 Paris, France

<sup>2</sup>Laboratoire de mathématiques d'Orsay, Université Paris-Saclay, CNRS, 91405 Orsay, France

<sup>3</sup>DataShape, Centre Inria Saclay, 91120 Palaiseau, France

December 7, 2020

## Abstract

The study of the topology of the superlevel sets of stochastic processes on  $[0, 1]$  in probability led to the introduction of  $\mathbb{R}$ -trees which characterize the connected components of the superlevel-sets. We provide a generalization of this construction to more general deterministic continuous functions on some path-connected, compact topological set  $X$  and reconcile the probabilistic approach with the traditional methods of persistent homology. We provide an algorithm which functorially links the tree  $T_f$  associated to a function  $f$  and study some invariants of these trees, which in 1D turn out to be linked to the regularity of  $f$ .

## 1 Introduction

A problem of interest in topological data analysis (TDA) is the study of the barcode of random functions. The easiest case-study of this problem is to look at some stochastic process  $S : [0, 1] \rightarrow \mathbb{R}$ , and study the barcodes of its sample paths. Studying the level sets and superlevel sets of stochastic processes in this particular setting has been widely studied by probabilists from different approaches [3, 4, 13, 15, 16, 18, 24], although some headway in understanding the allure of the barcode of Brownian motion has also been made by the TDA community [10]. When looking at the probability literature, one very quickly encounters trees which can be constructed from a deterministic function  $f : [0, 1] \rightarrow \mathbb{R}$ , and which exactly describe the topology of the (super)level sets of the function. This construction was originally thought of by Le Gall [15], who was able to describe fine properties of the sample paths of Lévy stochastic processes in [16] using this formalism.

---

<sup>\*</sup>Email: [daniel.perez@ens.fr](mailto:daniel.perez@ens.fr)

This work aims at reconciling the tree approach taken by the probability theory community with objects relevant to TDA, *i.e.* persistent diagrams or barcodes. Trees and barcodes turn out to be closely related. In fact, the link between the two is functorial. Establishing this connection has the merit to – among others – allow us to use the work of the probability theory community to further our understanding of the allure of barcodes of random processes. More importantly, it enlarges the TDA toolbox by looking at persistence diagrams as being themselves metric objects (trees), which yields new invariants for barcodes which are useful in a  $C^0$ -setting, such as the upper-box dimension.

## 1.1 State of the art

These trees and their construction have been extensively studied over  $[0, 1]$  [13, 15, 16, 24]. In particular, there are two results to keep in mind for the rest of this paper:

1. There is a link between the metric invariants of trees and the regularity of the function  $f$  from which they stem (theorem 3.7), this result was established by Picard for functions  $f : [0, 1] \rightarrow \mathbb{R}$  [24].
2. The  $L^\infty$ -stability of trees, with respect to the natural distance on the space of compact trees : the Gromov-Hausdorff distance (theorem 4.5). This is a result of Le Gall for functions on  $[0, 1]$  [15].

## 1.2 Our contribution

Our contribution can be summed up along the following four points.

1. A construction of trees which correctly computes the persistent homology of a function  $f : X \rightarrow \mathbb{R}$ , for any compact, connected and locally path connected topological space  $X$  (section 2.1). This is the natural generalization of the construction done by Le Gall on  $[0, 1]$ .
2. The explicit functorial link between the  $H_0$ -barcode of the function  $f$  and the tree constructed from this function,  $T_f$ , and do so constructively by virtue of algorithm 1.
3. An extension of the realm of validity of the known theorems for trees on  $[0, 1]$ , while providing links – whenever possible – to commonly studied quantities in persistent homology, such as the functional  $\ell_p$  (denoted  $\text{Pers}_p$  in the TDA literature) of [2, 8, 14, 29]. These extensions are theorems 3.7, 3.11 and 4.6.
4. A constructive proof of existence for the inverse problem, that is : given some tree  $T$  (of finite upper-box dimension), can we construct a function  $f : [0, 1] \rightarrow \mathbb{R}$  such that  $T_f = T$  ?

### 1.3 Layout of the paper

First, we construct the tree  $T_f$  associated to a continuous function  $f : X \rightarrow \mathbb{R}$ , on a connected and locally path-connected, compact topological space  $X$  by introducing a pseudo-distance  $d_f$  on  $X$  and defining  $T_f := X/\{d_f = 0\}$ . We will then discuss the relation between barcodes and trees and give the functorial relation between the two. Then, we give the link between the regularity of  $f : X \rightarrow \mathbb{R}$  and the metric properties of  $T_f$ . Subsequently, we provide a generalization of Le Gall's theorem on the stability of trees by relating  $d_{GH}(T_f, T_g)$  with  $\|f - g\|_{L^\infty}$ . Finally, we answer the question of the inverse problem positively and provide the construction of a function  $f : [0, 1] \rightarrow \mathbb{R}$  given a tree of finite upper-box dimension  $T$ . For the rest of the paper, we take all homology groups with respect to the field  $\mathbb{Z}/2\mathbb{Z}$ .

## 2 Trees and barcodes

### 2.1 Constructing a tree from a continuous function

For the rest of this paper, let  $X$  denote a path-connected, compact topological space and let  $f : X \rightarrow \mathbb{R}$  be a continuous function. Let us denote  $(X_r)_{r \in \mathbb{R}}$  the filtration of  $X$  by the **superlevels** of  $f$ , that is

$$X_r := \{x \in X \mid f(x) \geq r\}. \quad (2.1)$$

*Remark 2.1.* Notice that  $X_r$  is closed. In particular, despite  $X$  being locally path connected and connected,  $X_r$  might not be locally path connected. However, the interior of  $X_{r-\varepsilon}$  is for all  $\varepsilon > 0$  and  $X_r \subset \text{Int}(X_{r-\varepsilon})$ . It is in this sense that the arguments of this paper are to be understood whenever we use the local path connectedness of a superlevel  $X_r$ .

There exists a pseudo-distance on  $X$ , denoted  $d_f$ , given by:

**Definition 2.2.** Let  $X$  and  $f$  be defined as above. The  $H_0$ -**distance**,  $d_f$ , is the pseudo-distance

$$d_f(x, y) := f(x) + f(y) - 2 \sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f(\gamma(t)), \quad (2.2)$$

where the supremum runs over every path  $\gamma$  linking  $x$  to  $y$ .

*Remark 2.3.* Notice there are different ways of writing this distance. In particular, the sup above is also characterized by

$$\sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f(\gamma(t)) = \sup\{r \mid [x]_{H_0(X_r)} = [y]_{H_0(X_r)}\} \quad (2.3)$$

$$= \sup\{r \mid \exists \gamma \in C_1(X_r) \text{ such that } \partial \gamma = x - y\}. \quad (2.4)$$

These equalities hold, since we take the coefficients of homology with respect to  $\mathbb{Z}/2\mathbb{Z}$ , so we can interpret 1-cycles as sums of paths on  $X$ .

This pseudo-distance is a generalization of the distance introduced by Curien, Le Gall and Miermont in [13]. Note that  $d_f$  has the following properties:

1. **Identification of the connected components of superlevel sets:**  $d_f(x, y) = 0$  if and only if there exists  $t \in \mathbb{R}$  such that  $x, y \in \{f = t\}$  and  $[x]_{H_0(X_t)} = [y]_{H_0(X_t)}$ . In words, that if two points  $x$  and  $y$  are on the same level-set of  $f$  and they both lie in the same connected component of the super-level set above the level-set, then they are a distance zero away from one another.
2. **Compatibility with the filtration induced by  $f$ :** Let  $x, y \in X$  and suppose that  $f(x) < f(y)$ , then if  $[x]_{H_0(X_{f(x)})} = [y]_{H_0(X_{f(x)})}$ ,

$$d_f(x, y) := |f(x) - f(y)|. \quad (2.5)$$

Let us now prove that  $d_f$  indeed satisfies all of the statements above.

**Proposition 2.4.** The function  $d_f : X^2 \rightarrow \mathbb{R}^+$  of definition 2.2 is indeed a pseudo-distance.

*Proof.* Checking symmetry and positivity is easy. The only non-obvious point is that the triangle inequality is satisfied by this expression, so we will focus on this point. Let  $x, y, z \in X$  and denote

$$[x \mapsto y] := \sup_{\gamma: x \mapsto y} \inf_{t \in [0,1]} f \circ \gamma(t). \quad (2.6)$$

We will now show the following inequality

$$[x \mapsto z] + [z \mapsto y] \leq [x \mapsto y] + f(z), \quad (2.7)$$

which implies the triangle inequality. If we denote  $\gamma$  a path from  $x$  to  $z$  and  $\eta$  a path from  $z$  to  $y$  and by  $\gamma * \eta$  the concatenation of these two path, by definition, we have that

$$\inf_{t \in [0,1]} f \circ (\gamma * \eta)(t) \leq [x \mapsto y]. \quad (2.8)$$

Since this holds for any such two paths  $\gamma$  and  $\eta$ , it follows that

$$[x \mapsto z] \wedge [z \mapsto y] \leq [x \mapsto y]. \quad (2.9)$$

Without loss of generality, suppose that  $[x \mapsto z]$  achieves the above minimum. Furthermore, note that

$$[z \mapsto y] \leq f(z) \quad (2.10)$$

by definition of  $[z \mapsto y]$ . Adding the two last inequalities together we get that:

$$[x \mapsto z] + [z \mapsto y] \leq [x \mapsto y] + f(z), \quad (2.11)$$

as desired. ■

The compactibility with the filtration induced by  $f$  is immediate from the definition of  $d_f$ , so let us verify that  $d_f$  correctly identifies the connected components of superlevel sets.

**Proposition 2.5.** Let  $f$  be a continuous function as above, then  $d_f$  identifies the connected components of the superlevel sets.

*Proof.* We must check that  $d_f(x, y) = 0$  if and only if there exists  $t \in \mathbb{R}$  such that  $x, y \in \{f = t\}$  and  $[x]_{H_0(X_t)} = [y]_{H_0(X_t)}$ . The  $(\Leftarrow)$  direction is immediate, so let us focus on  $(\Rightarrow)$ .

Suppose that  $d_f(x, y) = 0$  and that  $f(x) \neq f(y)$ , then,

$$\sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f(\gamma(t)) = \frac{f(x) + f(y)}{2} > f(x) \wedge f(y). \quad (2.12)$$

However,

$$\sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f(\gamma(t)) \leq f(x) \wedge f(y), \quad (2.13)$$

which leads to a contradiction, so  $f(x) = f(y)$ . The condition  $d_f(x, y) = 0$  becomes:

$$f(x) = \sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f(\gamma(t)). \quad (2.14)$$

This is only possible if a path achieving the supremum lies entirely above  $f(x)$ . In other words, there is a path linking  $x$  and  $y$  contained within a connected component of  $X_{f(x)}$ , which implies the second condition. If there is no path achieving this supremum, we can proceed by an approximation argument and conclude that for every  $\varepsilon > 0$  there is a path lying entirely above  $f(x) - \varepsilon$ , implying that :

$$x, y \in \bigcap_{\varepsilon > 0} X_{f(x) - \varepsilon} \quad (2.15)$$

but since all  $X_{f(x) - \varepsilon}$  are closed and so is  $X_{f(x)}$ , this is only possible if  $x, y \in X_{f(x)}$ . ■

With these technicalities out of the way, let us consider the metric space

$$(T_f, d_f) := (X / \{d_f = 0\}, d_f), \quad (2.16)$$

where the quotient denotes the quotient of  $X$  where we identify all points  $x$  and  $y$  on  $X$  which satisfy  $d_f(x, y) = 0$ . Slightly abusing the notation let  $d_f$  denote the distance induced on the quotient by the pseudo-distance  $d_f$  on  $X$ .

The metric structure of  $T_f$  is simple as  $T_f$  is an  $\mathbb{R}$ -tree. Recall that  $\mathbb{R}$ -trees are defined as follows :

**Definition 2.6** (Chiswell, [11]). An  $\mathbb{R}$ -tree  $(T, d)$  is a connected metric space such that any of the following equivalent conditions hold:

- $T$  is a geodesic connected metric space and there is no subset of  $T$  which is homeomorphic to the circle,  $\mathbb{S}_1$ ;
- $T$  is a geodesic connected metric space and the Gromov 4-point condition holds, *i.e.* :

$$\forall x, y, z, t \in T \quad d(x, y) + d(z, t) \leq \max [d(x, z) + d(y, t), d(x, t) + d(y, z)] ;$$

- $T$  is a geodesic connected 0-hyperbolic space.

A **rooted  $\mathbb{R}$ -tree**  $(T, O, d)$  is an  $\mathbb{R}$ -tree along with a marked point  $O \in T$ .

Note that  $T_f$  is clearly connected, since  $X$  is connected. Our proof strategy will be to use the first characterization above and to split the proof into two parts. First, we will show that there are no subspaces of  $T_f$  which are homeomorphic to  $\mathbb{S}_1$ , and then that  $T_f$  is in fact a geodesic metric space.

Before proving this, it is helpful to introduce some notation for some maps and quantities which will become useful later on. Let  $\pi_f : X \rightarrow T_f$  denote the canonical projection onto  $T_f$  and let  $O$  denote the root of  $T_f$  (*i.e.*  $f(O) = \min f$ ), let us define the following quantities:

$$\ell(\tau) := \inf_X f + d_f(O, \tau) \quad (2.17)$$

$$h(\tau) := \sup_{x \in X_{f(\tau)}^\tau} f(x) - \ell(\tau) \quad (2.18)$$

where  $X_{f(\tau)}^\tau$  denotes the connected component of the superlevel set  $X_{f(\tau)}$  containing a preimage of  $\tau$ . These quantities are well-defined, by definition of  $d_f$ .

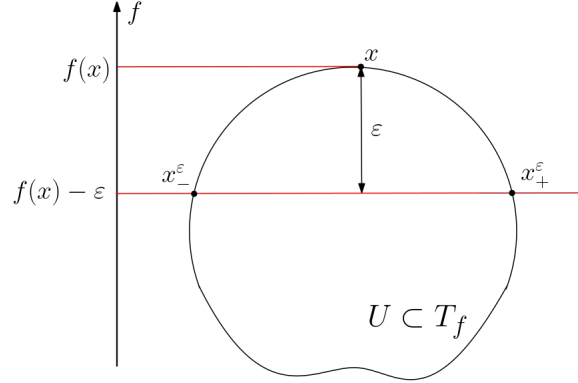
**Proposition 2.7.** The metric space  $T_f := X/\{d_f = 0\}$  equipped with distance  $d_f$  possesses no subspace homeomorphic to  $\mathbb{S}_1$ .

*Proof.* We will reason by contradiction. Suppose that  $T_f$  contains  $U \subset T_f$  such that  $U$  is homeomorphic to the circle,  $\mathbb{S}_1$ . Note that  $f$  descends to a function on  $T_f$  which is not locally constant anywhere by definition of  $d_f$ .

The restriction of  $f$  to  $U$  defines a function which admits some point  $x$  for which  $f$  is maximal on  $U$ . What we will now show is that, the filtration of the topological space  $X$  by  $f$  and the fact that  $d_f$  correctly identifies the connected components of superlevel sets of  $f$  forbid that we can “branch downwards” from  $x$ , thereby leading to a contradiction.

An element  $x \in T_f$  can be seen as a representative of a homology class  $[x]_{H_0(X_{f(x)})}$  lying at a certain level  $\{f = f(x)\}$ . For all  $\varepsilon > 0$  small enough, since  $U$  is homeomorphic to  $\mathbb{S}_1$  and since  $f$  is not locally constant on  $T_f$ , two different points  $x_+^\varepsilon$  and  $x_-^\varepsilon$  satisfying:

$$f(x_+^\varepsilon) = f(x_-^\varepsilon) = f(x) - \varepsilon \quad (2.19)$$



exist on  $U$ .

Notice that  $X_{f(x)-\varepsilon}^{x_+^\varepsilon}$  and  $X_{f(x)-\varepsilon}^{x_-^\varepsilon}$  are two connected components of  $X_{f(x)-\varepsilon}$ , so they are either equal or disjoint. If they are equal, then their image by  $\pi_f$  is the same, but  $x_+^\varepsilon$  and  $x_-^\varepsilon$  were supposed to be distinct in  $T_f$ , *i.e.* at a non-zero distance away from one another. But since all their preimages lie at the same level, and are in the same connected component of  $X_{f(x)-\varepsilon}$ ,  $x_+^\varepsilon$  and  $x_-^\varepsilon$  are in fact the same, leading to a contradiction. It follows that  $X_{f(x)-\varepsilon}^{x_+^\varepsilon}$  and  $X_{f(x)-\varepsilon}^{x_-^\varepsilon}$  must be disjoint.

If these two connected components are disjoint, then their image by  $\pi_f$  must be as well. Otherwise, there exists a point  $\tau$  on  $T_f$  lying in  $\pi_f(X_{f(x)-\varepsilon}^{x_+^\varepsilon})$  and  $\pi_f(X_{f(x)-\varepsilon}^{x_-^\varepsilon})$  simultaneously. However, every preimage of  $\tau$  in  $X$  lies on level  $f(\tau) > f(x) - \varepsilon$  and they are all in the same connected component of  $X_{f(\tau)}$ . The inclusion  $X_{f(\tau)} \hookrightarrow X_{f(x)-\varepsilon}$  is injective, so every preimage of  $\tau$  must lie in one and only one connected component of  $X_{f(x)-\varepsilon}$ . It follows that  $\pi_f(X_{f(x)-\varepsilon}^{x_+^\varepsilon})$  and  $\pi_f(X_{f(x)-\varepsilon}^{x_-^\varepsilon})$  are disjoint. Let us note:

$$T_r := \{\tau \in T_f \mid f(\tau) \geq r\} \quad (2.20)$$

These images are in fact nothing other than the connected component of  $T_{f(x)-\varepsilon}$  containing  $x_+^\varepsilon$  or  $x_-^\varepsilon$  respectively, as if  $z \notin X_{f(x)-\varepsilon}^{x_\pm^\varepsilon}$  then  $\pi_f(z) \notin T_{f(x)-\varepsilon}^{x_\pm^\varepsilon}$ . But, there is an arc linking  $x_+^\varepsilon$  and  $x_-^\varepsilon$  contained entirely in  $T_{f(x)-\varepsilon}$ , so these connected components are the same, leading to a contradiction. ■

**Proposition 2.8.** The metric space  $(T_f, d_f)$  is an  $\mathbb{R}$ -tree and we can choose a root for  $T_f$  to be the point  $O$  the image in  $T_f$  of a point  $x \in X$  for which the function  $f$  is minimal.

*Proof.* The only thing left to show is that  $T_f$  is a geodesic space. Let  $x$  and  $y$  be two points of  $X$ . If  $f(x) = f(y)$  and  $x$  and  $y$  are in the same connected component, there is nothing to show, so suppose that  $f(x) < f(y)$ . As before, note that  $f$  descends to the quotient and induces a non-locally constant function on  $T_f$ .



First, suppose that  $x$  and  $y$  are in the same connected component of  $X_{f(x)}$  and consider a path in  $X_{f(x)}$  going from  $y$  to  $x$ ,  $\gamma : [0, 1] \rightarrow X$ , which exists since  $X$  is path connected. The path  $\gamma$  can be modified into a path

$$\tilde{\gamma}(t) := \pi_f \left( \gamma \left( \arg \min_{s \in [0, t]} f \circ \gamma(s) \right) \right). \quad (2.21)$$

On this modified path  $f$  is decreasing implying that it does not self-intersect, although it may be locally constant. The length of  $\tilde{\gamma}$  is defined as

$$L(\tilde{\gamma}) = \sup_{(t_i)} \sum_{(t_i)} d_f(\tilde{\gamma}(t_{i+1}), \tilde{\gamma}(t_i)), \quad (2.22)$$

where the supremum is taken over all finite partitions of  $[0, 1]$ . For any finite partition, this sum is always bounded by  $f(y) - f(x)$ , since along  $\tilde{\gamma}$

$$f(\tilde{\gamma}(t_i)) \geq f(\tilde{\gamma}(t_{i+1})) \implies d_f(\tilde{\gamma}(t_{i+1}), \tilde{\gamma}(t_i)) = f(\tilde{\gamma}(t_i)) - f(\tilde{\gamma}(t_{i+1})) \quad (2.23)$$

by monotonicity of  $f$  along  $\tilde{\gamma}$ . This leads to pairwise cancelation of terms in the sum of equation 2.22. And so,

$$L(\tilde{\gamma}) = d_f(x, y). \quad (2.24)$$

Now, suppose that  $x$  and  $y$  are two points on  $X$ , such that  $f(x) \leq f(y)$  but such that  $x$  and  $y$  no longer lie in the same connected component of  $X_{f(x)}$  and pick a maximizer  $\gamma$  of the supremum (cf. remark 2.9)

$$\sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f \circ \gamma(t). \quad (2.25)$$

Since  $y$  is not connected to  $x$  in  $X_{f(x)}$ , by continuity of  $f$ , the path  $\gamma$  must eventually go under the level  $f(x)$ . Let us set

$$t^* := \sup \left\{ \arg \min_{s \in [0, 1]} f \circ \gamma(s) \right\} \quad (2.26)$$

and note that  $\gamma(t^*) < f(x)$ .

On  $[0, t^*]$ , the path  $\gamma$  lies entirely in  $X_{f(\gamma(t^*))}$  and similarly, entirely in  $X_{f(\gamma(t^*))}$  on  $[t^*, 1]$ . On  $[0, t^*]$ , we can define a modification of  $\gamma$ ,  $\tilde{\gamma} : [0, t^*] \rightarrow T_f$  by

$$\tilde{\gamma}(t) := \pi_f \left( \gamma \left( \arg \min_{s \in [0, t]} f \circ \gamma(s) \right) \right). \quad (2.27)$$

Analogously, if we define  $\eta(t) := \gamma(1 - t)$  – the reversed version of  $\gamma$  – it is possible to define a modification of  $\eta$ ,  $\tilde{\eta} : [0, 1 - t^*] \rightarrow T_f$ , by

$$\tilde{\eta}(s) := \pi_f \left( \eta \left( \arg \min_{r \in [0, s]} f \circ \eta(r) \right) \right). \quad (2.28)$$

In particular,  $\tilde{\eta}(1 - t^*) = \tilde{\gamma}(t^*)$ . If  $\tilde{\eta}_-$  denotes the reversed path along  $\tilde{\eta}$ , the concatenation (without reparametrization),

$$\zeta := \tilde{\gamma} * \tilde{\eta}_- \quad (2.29)$$

is a path going from  $\pi_f(x)$  to  $\pi_f(y)$  monotone decreasing on  $[0, t^*]$  and monotone increasing on  $]t^*, 1]$ .

For all  $\varepsilon > 0$ ,  $\zeta(t^* + \varepsilon)$  does not lie in the same connected component of  $X_{f(\zeta(t^* + \varepsilon))}$  as  $\pi_f(x)$ , but lies in the same connected component of  $X_{f(\zeta(t^* + \varepsilon))}$  as  $\pi_f(y)$ . We are thus reduced to examine the length of the path along two different sections of  $\zeta$ , each lying in the same connected component as either  $\pi_f(x)$  and  $\pi_f(y)$ . By the previous argument for points of  $T_f$  lying in the same connected component of a superlevel set, the length of  $\zeta$  is

$$L(\zeta) = f(x) - f(\zeta(t^*)) + f(y) - f(\zeta(t^*)) = d_f(x, y) \quad (2.30)$$

by definition of  $d_f(x, y)$ . Thus,  $T_f$  is indeed geodesic and it is an  $\mathbb{R}$ -tree, by virtue of proposition 2.7.

Finally, the tree is rooted since for any  $r < \inf f$ , every single point of  $X_r = X$  is identified in the quotient (since  $X$  was supposed to be connected), so we can identify the root with the point of  $T_f$  achieving this infimum. ■

*Remark 2.9.* If the suprema in the proof of proposition 2.8 are not achieved, it suffices to take a sequence of approximating paths  $(\gamma_n)$  and apply the same reasoning as above. For all  $\varepsilon > 0$  these approximations will yield paths of length  $d_f(x, y) + \varepsilon$ , from which the statement follows.

*Remark 2.10.* If  $X = [0, 1]$ , there is only one possible path between any two points  $x$  and  $y$ , so the definition above boils down to

$$d_f(x, y) := f(x) + f(y) - 2 \inf_{t \in [0, 1]} f, \quad (2.31)$$

which is exactly the distance originally introduced by Le Gall *et al.* [15].

## 2.2 An algorithm linking barcodes and trees

Given a tree stemming from a continuous function  $f : X \rightarrow \mathbb{R}$ , it is possible to reconstruct the  $H_0$ -barcode of  $f$  by only using  $T_f$ . If  $T_f$  is finite, the relation between the barcode of  $H_0(X, f)$  with respect to the superlevel filtration and the tree  $T_f$  is given by algorithm 1.

If  $T_f$  is infinite, we can still give a correspondence between the barcode and the tree proceeding by approximation. This approximation procedure requires the introduction of so-called  $\varepsilon$ -trimmings of  $T_f$ , of which we briefly recall the definition.

**Definition 2.11.** The  $\varepsilon$ -simplified tree of  $f$ ,  $T_f^\varepsilon$  or the  $\varepsilon$ -trimmed tree of  $f$ , is the subtree of  $T_f$  defined as

$$T_f^\varepsilon := \{\tau \in T_f \mid h(\tau) \geq \varepsilon\} \quad (2.32)$$

---

**Algorithm 1:** A functorial relation between persistent modules and  $\mathbb{R}$ -trees

---

**Result:**  $\mathbb{V}$   
 $\mathcal{F} \leftarrow T$  ;  
 $\mathbb{V} \leftarrow 0$  ;  
 $i \leftarrow 0$  ;  
**while**  $\mathcal{F} \neq \emptyset$  **do**  
    Find  $\gamma$  the longest path in  $\mathcal{F}$  starting from a root  $\alpha$  and ending in a leaf  $\beta$  ;  
    **if**  $i = 0$  **then**  
         $\mathbb{V} \leftarrow \mathbb{V} \oplus k[\ell(\alpha), \infty[$  ;  
    **else**  
         $\mathbb{V} \leftarrow \mathbb{V} \oplus k[\ell(\alpha), \ell(\beta)[$  ;  
    **end**  
     $\mathcal{F} \leftarrow \overline{\mathcal{F} \setminus \text{Im}(\gamma)}$  ;  
     $i \leftarrow i + 1$  ;  
**end**  
**return**  $\mathbb{V}$

---

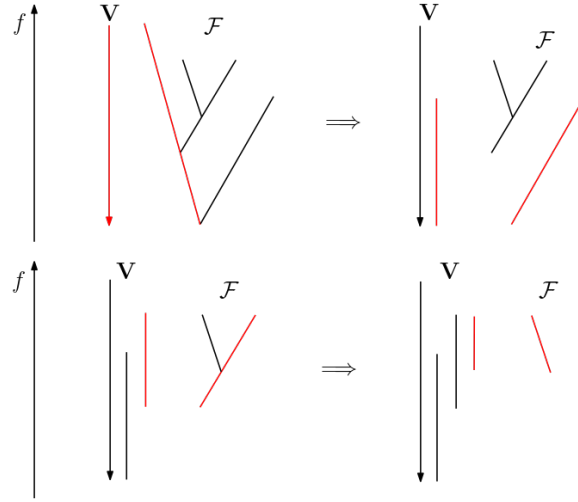


Figure 1: The first four iterations of algorithm 1. For every step, in red is the longest branch of the tree, which we use to progressively construct the persistent module  $\mathbb{V}$  by associating an interval module whose ends correspond exactly to the values of the endpoints of the branches.

*Remark 2.12.* The definition above can be extended to hold for any rooted tree  $(T, d, O)$  if we give define an equivalent for  $h(\tau)$ . This amounts to defining a suitable filtration of  $T$  by some function  $f$ . For any tree, we will take this function  $f$  to be  $f(\tau) = d(O, \tau)$ .

An  $\varepsilon$ -trimmed tree is always finite by virtue of the compactness of  $X$ . For a monotone decreasing sequence  $(\varepsilon_n)_{n \in \mathbb{N}}$  such that  $\varepsilon_n \rightarrow 0$ , we have the following chain of inclusions

$$T^{\varepsilon_1} \hookrightarrow T^{\varepsilon_2} \hookrightarrow T^{\varepsilon_3} \hookrightarrow \dots \quad (2.33)$$

Applying algorithm 1, we get a set of maps on the persistence modules induced by these inclusions. More precisely, denoting  $\text{Alg}(T_f^{\varepsilon_n})$  the output of the algorithm

$$\text{Alg}(T^{\varepsilon_1}) \rightarrow \text{Alg}(T^{\varepsilon_2}) \rightarrow \text{Alg}(T^{\varepsilon_3}) \rightarrow \dots \quad (2.34)$$

where the morphisms are the maps induced at the level of the interval modules generating  $\text{Alg}(T^{\varepsilon_n})$ . Indeed, the interval modules  $k[\alpha, \beta_n[$  of  $\text{Alg}(T^{\varepsilon_n})$  satisfy that there is exactly one interval module of  $\text{Alg}(T^{\varepsilon_m})$  ( $m > n$ ) such that  $[\alpha, \beta_n[ \subset [\alpha, \beta_m[$ . A natural definition for infinite  $T$  is thus

$$\text{Alg}(T) := \varinjlim \text{Alg}(T_f^{\varepsilon_n}) \quad (2.35)$$

In categoric terms, the algorithm above in fact is a functor

$$\text{Alg} : \mathbf{Tree} \rightarrow \mathbf{PersMod}_k, \quad (2.36)$$

where  $\mathbf{Tree}$  is the category of rooted  $\mathbb{R}$ -trees seen as metric spaces, whose morphisms are isometric embeddings (which are not required to be surjective) preserving the roots, and where  $\mathbf{PersMod}_k$  is the category of q-tame persistence modules over a field  $k$  (cf. Oudot's book for details on the category of persistence modules [23]). The action of  $\text{Alg}$  on morphisms between two trees  $\zeta : T \rightarrow T'$  is defined as follows. If both  $T$  and  $T'$  are finite, since  $\zeta$  is an isometric embedding and it preserves the root, we can define  $\text{Alg}(\zeta)$  to be

$$\text{Alg}(\zeta) := \bigoplus_i \text{id}_{k[\zeta(\alpha_i), \zeta(\beta_i)[}, \quad (2.37)$$

where  $k[\alpha_i, \beta_i[$  denotes the modules in the interval module decomposition of  $\text{Alg}(T)$  (which is finite, since  $T$  is as well). If  $T$  is infinite, we extend the above definition by taking successive  $\varepsilon_n$ -simplifications of  $T$  and taking the direct limit of the construction above. Note that this procedure is well-defined since  $\varepsilon_n$ -simplifications only depend on the function  $h$ , which in turn can be taken to only depend on the distance to the root.

*Remark 2.13.* To define  $\text{Alg}$ , we do not need the tree  $T$  to stem from a function  $f$ , since the algorithm only depends on the function  $\ell$ , which we can define to be the distance from the root to a point  $\tau \in T$ .

Let us now consider a tree  $T_f$  stemming from a function  $f$  and show that  $\text{Alg}(T_f) = H_0(X, f)$ . To do this, we will need the following proposition.

**Proposition 2.14.** Let  $\tau$  and  $\eta$  be elements of  $T_f$  such that  $f(\tau) < f(\eta)$  and let  $x \in \pi^{-1}(\tau)$  and  $y \in \pi^{-1}(\eta)$ , then

$$\exists \text{ path } \gamma : x \mapsto y \text{ s.t. } \forall t, f(\gamma(t)) \geq f(\tau) \iff h(\tau) \geq f(\eta) - f(\tau) \text{ and } x, y \in X_{f(\tau)}^\tau. \quad (2.38)$$

*Proof.* Since there exists  $\gamma$  connecting  $x$  and  $y$  and since  $\gamma$  always stays above  $f(\tau)$ , we conclude naturally that  $\text{Im}(\gamma) \subset X_{f(\tau)}^\tau$ , which implies that  $h(\tau) \geq f(\eta) - f(\tau)$  by definition of  $h(\tau)$ .

The implication ( $\Leftarrow$ ) is clear since if  $x, y \in X_{f(\tau)}^\tau$  and  $X_{f(\tau)}^\tau$  is connected, by path connexity of  $X$  there exists a path between  $x$  and  $y$  which stays above  $f(\tau)$ . ■

This proposition suffices to prove the following theorem on the validity of algorithm 1.

**Theorem 2.15.** Let  $f : X \rightarrow \mathbb{R}$  be continuous. Then  $\text{Alg}(T_f) = H_0(X, f)$ .

*Proof.* Suppose that  $T_f$  is finite. If this is the case, then  $\text{Alg}(T_f)$  is a decomposable persistence module  $\text{Alg}(T_f) := \mathbb{V}$ . The fact that  $\mathbb{V}$  is pointwise isomorphic to  $H_0(X, f)$  holds since  $d_f$  correctly identifies the connected components of the superlevel sets. This guarantees the existence of a pointwise isomorphism since both spaces have the same (finite) dimension.

Let us now check that  $\text{rank}(\mathbb{V}(r \rightarrow s)) = \text{rank}(H_0(X_r \rightarrow X_s))$ . The inclusion  $X_r \hookrightarrow X_s$  induces the following long exact sequence in homology

$$\begin{array}{ccccccc} \cdots & \rightarrow & H_1(X_s) & \rightarrow & H_1(X_s, X_r) & \rightarrow & H_0(X_r) \\ & & & & \downarrow & & \\ & & & & H_0(X_s) & \rightarrow & H_0(X_s, X_r) \rightarrow 0 \end{array}$$

Since this sequence is exact

$$\text{rank}(H_0(X_r \rightarrow X_s)) = \dim \ker(H_0(X_s) \rightarrow H_0(X_s, X_r)). \quad (2.39)$$

For notational simplicity, let us denote  $\phi : H_0(X_s) \rightarrow H_0(X_s, X_r)$ . Note that  $\phi[c] = [0]$  if and only if there is a path  $\gamma$  between the representative  $c \in X_s$  and an element  $b \in X_r$  such that  $\gamma$  stays within  $X_s$ . Without loss of generality, let us take  $c$  such that  $c \in \{f = s\}$ . Finding such a path  $\gamma$  is only possible if  $c$  and  $b$  lie in the same connected component of  $X_r$ . By proposition 2.14, this can happen if and only if  $h([c]_{T_f}) \geq r - s$ . It follows that

$$\dim \ker \phi = \#\{\tau \in T_f \mid h(\tau) \geq r - s \text{ and } f(\tau) = s\}, \quad (2.40)$$

which concludes the proof for the finite case.

If  $T_f$  is infinite, we consider a sequence of  $\varepsilon_n$ -trimmings of  $T_f$  such that  $\varepsilon_n \xrightarrow{n \rightarrow \infty} 0$ . For any  $r > s$ , there exists  $n$  such that  $r - s > \varepsilon_n$ . But  $T_f^{\varepsilon_n}$  is finite, so we are reduced to the previous case. ■

### 3 Regularity of $f$ and metric properties of $T_f$

From the general theory of persistence modules [9, 23], we already know that for all  $C^0$ -functions over a compact set  $X$ ,  $H_0(X, f)$  is a q-tame persistence module. The previous section is a particular case of this general fact in degree 0 in homology. Of course, if this was the only conclusion we could draw from the tree construction, the theory we have just exposed above would be completely redundant. However, the metric structure on  $T_f$  is richer than the persistence diagram structure of  $f$ . This is because the tree construction gives us access to a new family of (metric) invariants for barcodes.

In this section, we will give an interpretation of some of the metric invariants of  $T_f$ . While these results originally stemmed from probability theory [15, 24] they are completely deterministic. In 1D, the small bars of the barcode are related to the regularity of the function. The correct quantity to characterize this regularity in 1D is the true or total  $p$ -variation of the function  $f$ , of which we briefly recall the definition.

**Definition 3.1.** Let  $f : [0, 1] \rightarrow \mathbb{R}$  be a continuous function. The **true  $p$ -variation of  $f$**  is defined as

$$\|f\|_{p\text{-var}} := \left[ \sup_D \sum_{t_k \in D} |f(t_k) - f(t_{k-1})|^p \right]^{1/p}, \quad (3.41)$$

where the supremum is taken over all finite partitions  $D$  of the interval  $[0, 1]$ .

*Remark 3.2.* We speak about *true*  $p$ -variation to make the distinction with another notion of variation typically considered in the probabilistic context (more precisely, stochastic calculus), where instead of the supremum over all partitions, we have a limit as the mesh of the partition considered tends to zero.

This true  $p$ -variation has the nice property of being closely related to an  $\ell_p$  functional of the barcode, which has been used extensively by the TDA community where it is typically denoted  $\text{Pers}_p$  [8, 12, 14, 22, 29], it is defined as:

**Definition 3.3.** Let  $f : X \rightarrow \mathbb{R}$  be a continuous function. The  **$\ell_p$ -length of the barcode  $\mathcal{B}(f)$**  is the following functional

$$\ell_p(f) := \left( \sum_{I \in \mathcal{B}(f)} \mu(I \cap [\inf(f), \sup(f)])^p \right)^{1/p}, \quad (3.42)$$

where  $\mu$  is the Lebesgue measure on  $\mathbb{R}$  and where  $\mathcal{B}$  denotes the barcode of  $f$  obtained from the superlevel filtration.

It turns out that the  $\ell_p$  bounds the total  $p$ -variation of the function as can be shown by adapting and slightly tweaking a result by Picard [24]:

**Theorem 3.4** (Picard, §3 [24]). *Let  $f : [0, 1] \rightarrow \mathbb{R}$  be a continuous function, then  $\|f\|_{p-var}$  is finite as soon as  $\ell_p(f)$  is finite. In fact, for any  $p$*

$$\|f\|_{p-var}^p \leq 2\ell_p(f)^p. \quad (3.43)$$

*Furthermore, if  $\|f\|_{(p-\delta)-var}$  is finite for some  $\delta > 0$ ,  $\ell_p(f)$  is also finite.*

Intuition would have us think that generalizing this sort of statement relating regularity to the  $\ell_p$ -length of  $\mathcal{B}(f)$  beyond dimension 1 would be straightforward, since this norm is sensitive to the oscillations of the function. This turns out not to be the case, and in fact remains an open problem. One of the main difficulties is that it is not so clear how we should define an equivalent of the total variation on a more general topological space  $X$ , which is compatible enough with the study of superlevel sets to be useful.

For smooth  $f$ , there also seems to be a close link between the functionals  $\ell_p$  and different notions of regularity. For instance, on  $\mathbb{S}_1$ ,  $\ell_1$  is the total variation of  $f$  and on  $\mathbb{T}^2$ , Polterovitch *et al.* have given a result [26] relating  $\ell_1$  with the Sobolev  $W^{2,2}$ -norm. Unfortunately, the proof of this result heavily relies on the geometry of the torus and other properties which hold only in 2D.

**Theorem 3.5** (Polterovich, *et al.* [26]). *For every function  $f : \mathbb{T}^2 \rightarrow \mathbb{R}$  :*

$$\ell_1(f) \leq C \|f\|_{W^{2,2}}, \quad (3.44)$$

*where  $\|\cdot\|_{W^{2,2}}$  denotes the Sobolev  $(2, 2)$ -norm.*

These links between the functionals  $\ell_p$  and the regularity of  $f$  motivate the study of this functional in more detail. This is also suggested by the following result by Picard [24]:

**Theorem 3.6** (Picard, §3 [24]). *Let  $f : [0, 1] \rightarrow \mathbb{R}$  be a continuous function and denote*

$$\mathcal{V}(f) := \inf\{p \mid \|f\|_{p-var} < \infty\} \quad \text{and} \quad \mathcal{L}(f) := \inf\{p \mid \ell_p(f) < \infty\}. \quad (3.45)$$

*Then,*

$$\mathcal{V}(f) = \mathcal{L}(f) = \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)} = \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 = \overline{\dim} T_f \quad (3.46)$$

*where  $a \vee b := \max\{a, b\}$ ,  $N^\varepsilon$  is the number of leaves of the  $\varepsilon$ -trimmed tree  $T_f^\varepsilon$ ,  $\lambda(T_f^\varepsilon)$  denotes the length of  $T_f^\varepsilon$  and  $\overline{\dim}$  denotes the upper-box dimension.*

This result generalizes for more general topological spaces  $X$  and provides a connection between the topology of  $T_f$  and the  $\ell_p$  functional.

**Theorem 3.7.** *With the same notation as above and supposing that  $\overline{\dim} T_f$  is finite, the following chain of equalities hold*

$$\mathcal{L}(f) = \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 = \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)} = \overline{\dim} T_f. \quad (3.47)$$

Furthermore,

$$\underline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 \leq \underline{\dim} T_f \leq \underline{\lim}_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)}, \quad (3.48)$$

where  $\underline{\dim}$  is the lower-box dimension. For  $\underline{\dim} T_f > 1$ , these inequalities turn into equalities if either:

$$\overline{\lim}_{\varepsilon \rightarrow 0} \frac{N^{2\varepsilon}}{N^\varepsilon} < 1 \quad \text{or} \quad \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\lambda(T_f^{2\varepsilon})}{\lambda(T_f^\varepsilon)} < 1. \quad (3.49)$$

*Remark 3.8.* The study of  $N^\varepsilon$  is in fact completely equivalent to the study of  $\ell_p^p(f)$ . Indeed,

$$\ell_p^p(f) = p \int_0^\infty \varepsilon^{p-1} N^\varepsilon d\varepsilon, \quad (3.50)$$

which is finite as soon as  $p > \mathcal{L}(f)$ . This is nothing other than the Mellin transform of  $N^\varepsilon$ . By the Mellin inversion theorem, for any  $c > \mathcal{L}(f)$ , we have

$$N^\varepsilon = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \ell_p^p(f) \varepsilon^{-p} \frac{dp}{p}. \quad (3.51)$$

*Proof of theorem 3.7.* By the procedure detailed in section 5, since  $\overline{\dim} T_f$  is finite we can construct a function  $\hat{f} : [0, 1] \rightarrow \mathbb{R}$  such that  $T_f$  and  $T_{\hat{f}}$  are isometric. Applying Picard's theorem to  $T_{\hat{f}}$  and noting that  $\mathcal{L}(f)$  depends only on the  $T_f$ , we have that

$$\mathcal{L}(f) = \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)} = \overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 = \overline{\dim} T_f. \quad (3.52)$$

Let us now show the inequalities for  $\underline{\lim}$ . First,

$$\lambda(T_f^\varepsilon) = \int_\varepsilon^\infty N^a da, \quad (3.53)$$

which implies that

$$\underline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 \leq \underline{\lim}_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)}. \quad (3.54)$$

Additionally,

$$N^\varepsilon \leq \mathcal{N}(\varepsilon/2) \quad (3.55)$$



where  $\mathcal{N}(\varepsilon)$  denotes the minimal number of balls of radius  $\varepsilon$  necessary to cover  $T_f$ . This inequality holds as above each leaf of  $T_f^\varepsilon$ , at least one ball of radius  $\frac{\varepsilon}{2}$  is necessary to cover this section of the tree. It follows that

$$\lim_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 \leq \underline{\dim} T_f. \quad (3.56)$$

We can bound this minimal number of balls  $\mathcal{N}(\varepsilon)$  by the following

$$\mathcal{N}(\varepsilon) \leq N^{\varepsilon/2} + \frac{\lambda(T_f^{\varepsilon/2})}{\varepsilon/2} \leq 2 N^{\varepsilon/2} \vee \left\lceil \frac{\lambda(T_f^{\varepsilon/2})}{\varepsilon/2} \right\rceil, \quad (3.57)$$

which holds since, at most  $N^\varepsilon$  balls are needed to cover  $T_f \setminus T_f^\varepsilon$ . To cover  $T_f^\varepsilon$ , at most:  $\left\lceil \lambda(T_f^{\varepsilon/2})/(\varepsilon/2) \right\rceil$  balls are needed, so the inequality above follows by further majoring the terms. This implies that

$$\underline{\dim} T_f \leq \left[ \lim_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 \right] \vee \left[ \lim_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)} \right], \quad (3.58)$$

but by inequality 3.56 this means that

$$\underline{\dim} T_f \leq \lim_{\varepsilon \rightarrow 0} \frac{\log(\lambda(T_f^\varepsilon)/\varepsilon)}{\log(1/\varepsilon)}. \quad (3.59)$$

Finally,

$$\begin{aligned} \frac{\lambda(T_f^\varepsilon) - \lambda(T_f^{2\varepsilon})}{\varepsilon} &= \frac{1}{\varepsilon} \left[ \int_\varepsilon^\infty N^a da - \int_{2\varepsilon}^\infty N^a da \right] \\ &= \frac{1}{\varepsilon} \int_\varepsilon^{2\varepsilon} N^a da \leq N^\varepsilon, \end{aligned} \quad (3.60)$$

since  $N^\varepsilon$  is monotone decreasing. This reasoning also gives a lower bound:

$$N^{2\varepsilon} \leq \frac{\lambda(T_f^\varepsilon) - \lambda(T_f^{2\varepsilon})}{\varepsilon} \leq N^\varepsilon. \quad (3.61)$$

which entails that:

$$\lim_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} = \lim_{\varepsilon \rightarrow 0} \frac{\log \left[ \frac{\lambda(T_f^\varepsilon) - \lambda(T_f^{2\varepsilon})}{\varepsilon} \right]}{\log(1/\varepsilon)}. \quad (3.62)$$

Suppose that this limit is larger than 1. Rearranging, we get:

$$\frac{\varepsilon N^{2\varepsilon}}{\lambda(T_f^\varepsilon)} \leq 1 - \frac{\lambda(T_f^{2\varepsilon})}{\lambda(T_f^\varepsilon)} \leq \frac{\varepsilon N^\varepsilon}{\lambda(T_f^\varepsilon)}. \quad (3.63)$$

It follows that if any of these quantities admits a  $\underline{\lim}$  which is strictly greater than zero, we have

$$\underline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} = \underline{\lim}_{\varepsilon \rightarrow 0} \frac{\log \lambda(T_f^\varepsilon)}{\log(1/\varepsilon)}. \quad (3.64)$$

Another equivalent condition for the validity of this equality is whether

$$\underline{\lim}_{\varepsilon \rightarrow 0} \frac{N^{2\varepsilon}}{N^\varepsilon} < 1, \quad (3.65)$$

which finishes the proof.  $\blacksquare$

*Remark 3.9.* If  $\overline{\dim} = \underline{\dim}$ , all the limits of the above theorem are well-defined, yielding exact asymptotics for  $\lambda(T_f^\varepsilon)$  and  $N^\varepsilon$ . This is in particular the case if  $\overline{\dim} = \dim_H$ , where  $\dim_H$  denotes the Hausdorff dimension.

The functional  $\lambda(T_f^\varepsilon)$  is what some authors [25, 26] refer to as the Banach indicatrix and its asymptotics have a topological interpretation as described in the statement of the theorem. It is interesting to note that the study of the upper-box dimension is natural in the tree approach. Nonetheless,  $\underline{\dim}$  has also been used in the context of persistent homology in degree 0, albeit under a completely different setting by Schweinhart *et al.* in [20, 27] and by Adams *et al.* [1].

### 3.1 Another extension of Picard's theorem

It is possible to further extend Picard's theorem by some rudimentary considerations and by imposing the so-called locally linearly connected condition  $X$ . Let us briefly recall its definition.

**Definition 3.10.** A **locally linearly connected (LLC) metric space**  $(X, d)$ , is a metric space such that for all  $r > 0$  and for all  $z \in X$ , for all  $x, y \in B(z, r)$ , there exists an arc connecting  $x$  and  $y$  such that the diameter of this arc is linear in  $d(x, y)$ .

With this extra assumption, it is possible to link the regularity of the function  $f$  to the quantities defined in theorem 3.7.

**Theorem 3.11** (The regularity of  $f$  bounds the upper-box dimension of  $T_f$ ). *Let  $X$  be a compact LLC metric space. Keeping the same notations as in theorem 3.7, the following inequality holds*

$$\overline{\dim} T_f \leq \mathcal{H}(f) \overline{\dim} X, \quad (3.66)$$

where:

$$\mathcal{H}(f) := \inf \left\{ \frac{1}{\alpha} \mid \|f\|_{C^\alpha} < \infty \right\} \quad (3.67)$$

The proof of this theorem relies on two lemmas:

**Lemma 3.12.** Let  $X$  and  $Y$  be two metric spaces such that there is a surjective map  $\pi : X \rightarrow Y$  such that  $\pi \in C^\alpha(X, Y)$ , then

$$\overline{\dim} Y \leq \frac{1}{\alpha} \overline{\dim} X. \quad (3.68)$$

**Lemma 3.13.** Let  $X$  be a compact locally linearly connected (LLC) metric space (cf. definition 3.10) and let  $f : X \rightarrow \mathbb{R}$  be a continuous function, then

$$f \in C^\alpha(X, \mathbb{R}) \implies \pi_f \in C^\alpha(X, T_f). \quad (3.69)$$

Let us show that these two lemmas imply the theorem.

*Proof of theorem 3.11.* If  $f \notin C^\alpha(X, \mathbb{R})$  for any  $\alpha$ , there is nothing to show, since the statement is vacuous. Otherwise, the projection onto the tree of  $f$ ,  $\pi_f : X \rightarrow T_f$  is in  $C^\alpha(X, T_f)$  according to lemma 3.13. It follows from lemma 3.12 that

$$\overline{\dim} T_f \leq \frac{1}{\alpha} \overline{\dim} X. \quad (3.70)$$

The statement of the theorem follows by taking the infimum over all such  $\alpha$ . ■

All that remains to show is the two lemmas.

*Proof of lemma 3.12.* Since  $\pi : X \rightarrow Y$  is surjective and  $C^\alpha(X, Y)$ , for any  $x \in X$

$$\pi \left( B_X \left( x, \left( \frac{\varepsilon}{K} \right)^{1/\alpha} \right) \right) \subset B_Y(\pi(x), \varepsilon) \quad (3.71)$$

for some constant  $K$ . It follows that the minimal number of balls needed to cover  $X$ ,  $\mathcal{N}_X$  dominates the minimal number of balls needed to cover  $Y$ ,  $\mathcal{N}_Y$ . More precisely

$$\mathcal{N}_Y(\varepsilon) \leq \mathcal{N}_X \left( \left( \frac{\varepsilon}{K} \right)^{1/\alpha} \right) \iff \alpha \frac{\mathcal{N}_Y(\varepsilon)}{\log(1/\varepsilon) + \log(K)} \leq \frac{\mathcal{N}_X \left( \left( \frac{\varepsilon}{K} \right)^{1/\alpha} \right)}{\log \left( \left( \frac{K}{\varepsilon} \right)^{1/\alpha} \right)}.$$

The statement of the lemma follows. ■

*Proof of lemma 3.13.* Suppose that  $f : X \rightarrow \mathbb{R}$  is in  $C^\alpha(X, \mathbb{R})$  and let  $x, y \in X$  be two points inside a ball of radius  $r > 0$ . Without loss of generality, suppose that  $f(x) < f(y)$ . Since  $T_f$  is a geodesic space, the distance  $d_f(\pi_f(x), \pi_f(y))$  is the length of the geodesic arc in  $T_f$  linking  $\pi_f(x)$  and  $\pi_f(y)$ . By compactness of this geodesic path, there is a point  $\tau \in T_f$  where  $f$  achieves its minimum, thus

$$d_f(\pi_f(x), \pi_f(y)) = |f(y) - f(x)| + 2|f(x) - f(\tau)|. \quad (3.72)$$

This minimum  $f(\tau)$  has the particularity that

$$f(\tau) = \sup_{\gamma: x \rightarrow y} \inf_{t \in [0,1]} f \circ \gamma, \quad (3.73)$$

where the supremum is taken over all paths on  $X$  linking  $x$  and  $y$ . The equality holds, since it is always possible to find a path in  $X_r$  whose image in  $T_f$  contains a path in  $T_f$  lying entirely above level  $r$ . Since  $X$  is LLC, there is a path  $\eta$  linking  $x$  and  $y$  in  $X$  whose diameter we can control linearly in terms of  $d_X(x, y)$ . Along  $\eta$ , we have that

$$f(\tau) = \sup_{\gamma: x \rightarrow y} \inf_{t \in [0,1]} f \circ \gamma \geq \inf_{t \in [0,1]} f \circ \eta =: f(z) \quad (3.74)$$

for some  $z \in X$  along the path  $\eta$ . Combining the LLC condition and the fact that  $f$  is  $\alpha$ -Hölder gives:

$$f(x) - f(\tau) \leq f(x) - f(z) \leq d(x, z)^\alpha \leq C d(x, y)^\alpha \quad (3.75)$$

for some constant  $C$ , which is determined by quantitative LLC condition. Putting everything together we have that:

$$d_f(\pi_f(x), \pi_f(y)) \leq (2C + 1) d_X(x, y)^\alpha, \quad (3.76)$$

which finishes the proof. ■

Theorem 3.11 is sharp. Indeed, Brownian sample paths almost surely saturate this inequality. However, there is no hope to prove equality in all generality. Indeed note that for any  $f \in C^1(\mathbb{T}^2, \mathbb{R})$   $T_f$  is a finite tree and has upper-box dimension 1, but :

$$\overline{\dim} T_f = 1 < 2 = \mathcal{H}(f) \overline{\dim} \mathbb{T}^2. \quad (3.77)$$

An interesting problem would be to either prove that for irregular functions equality holds, or alternatively, to find a counter example of an irregular function for which we don't have equality. Of course, this can only be done in dimensions  $\geq 2$ . More precisely, we can phrase this question in the form of the following conjecture.

**Conjecture 3.14.** Given an LLC space  $X$ , we have:

$$\overline{\dim} X = \sup_{f \in C^0(X, \mathbb{R})} \frac{\overline{\dim} T_f}{\mathcal{H}(f)}. \quad (3.78)$$

Furthermore, there may be instances of spaces  $X$  for which this supremum is achieved by some  $f \in C^\alpha(X, \mathbb{R})$  for some  $0 < \alpha \leq 1$ .

## 4 Stability of trees with respect to the $L^\infty$ -norm

It is well known that the Gromov-Hausdorff distance is a natural notion of distance between metric spaces. Recall that this distance is defined as:

**Definition 4.1.** Let  $X$  and  $Y$  be two compact metric spaces, the **Gromov-Hausdorff distance**,  $d_{GH}(X, Y)$  between  $X$  and  $Y$ , is defined as

$$d_{GH}(X, Y) := \inf_{\substack{f: X \rightarrow Z \\ g: Y \rightarrow Z}} \max \left\{ \sup_{x \in X} \inf_{y \in Y} d_Z(f(x), g(y)), \sup_{y \in Y} \inf_{x \in X} d_Z(f(x), g(y)) \right\}. \quad (4.79)$$

where the infimum is taken over all metric spaces  $Z$  and all isometric embeddings  $f : X \rightarrow Z$  and  $g : Y \rightarrow Z$ .

From this definition, we see that  $d_{GH}$  quantifies how far away two metric spaces  $X$  and  $Y$  are from being isometric to each other. However, this definition, while practical in theory is very difficult to compute in practice. To somewhat alleviate this, we will use the following characterization of the Gromov-Hausdorff distance:

**Proposition 4.2** (Burago *et al.*, §7 [7]). The Gromov-Hausdorff distance is characterized by the following equality:

$$d_{GH}(X, Y) = \frac{1}{2} \inf_{\mathfrak{R}} \sup_{\substack{(x, y) \in \mathfrak{R} \\ (x', y') \in \mathfrak{R}}} |d_X(x, x') - d_Y(y, y')|, \quad (4.80)$$

where the infimum is taken over all *correspondences*, *i.e.* subsets  $\mathfrak{R} \subset X \times Y$  such that for every  $x \in X$  there is at least one  $y \in Y$  such that  $(x, y) \in \mathfrak{R}$  and a symmetric condition for every  $y \in Y$ .

*Remark 4.3.* Given two surjective maps  $\pi_X : Z \rightarrow X$  and  $\pi_Y : Z \rightarrow Y$ , it is possible to build a correspondence between  $X$  and  $Y$  by considering the set  $\{(\pi_X(z), \pi_Y(z)) \in X \times Y \mid z \in Z\}$ .

Recall that there is also a natural notion of distance in the space of barcodes (or equivalently, persistence diagrams) called the bottleneck distance, which we will note  $d_b$ . For the definition of this distance, we refer the reader to the books by Chazal and Oudot on persistence theory [9, 23]. With respect to this distance, we have a “stability theorem” which we briefly state:

**Theorem 4.4** (Bottleneck stability with respect to  $L^\infty$ , Corollary 3.6 [23]). *Let  $f, g : X \rightarrow \mathbb{R}$  be two continuous functions, then*

$$d_b(\mathcal{B}(f), \mathcal{B}(g)) \leq \|f - g\|_{L^\infty} \quad (4.81)$$

where  $\mathcal{B}(f)$  and  $\mathcal{B}(g)$  denote the barcodes (or diagrams) of  $f$  and  $g$  respectively.

A natural question is to ask whether we have an equivalent statement about the stability of  $d_{GH}$  with respect to  $\|\cdot\|_{L^\infty}$  and whether the two notions of distances are in some sense “compatible”. We will positively answer this first question. However, in general  $d_b$  and  $d_{GH}$  are not compatible, in the sense that no inequality between the two holds in all generality (*cf.* remark 4.7). Le Gall and Duquesne [15] gave a first stability result of  $d_{GH}$  respect to the  $L^\infty$ -norm on continuous functions on  $[0, 1]$ :

**Theorem 4.5** ( $L^\infty$ -stability of trees, [15]). *Let  $f, g : [0, 1] \rightarrow \mathbb{R}$  be continuous functions. Then*

$$d_{GH}(T_f, T_g) \leq 2 \|f - g\|_{L^\infty} . \quad (4.82)$$

This result for functions on  $[0, 1]$  generalizes to more general topological spaces  $X$ .

**Theorem 4.6** (Stability theorem for trees). *Let  $f$  and  $g : X \rightarrow \mathbb{R}$  be continuous, then*

$$d_{GH}(T_f, T_g) \leq 2 \|f - g\|_{L^\infty} . \quad (4.83)$$

*Proof.* We will use the distortion characterization of the Gromov-Hausdorff distance, which yields the following inequality

$$d_{GH}(T_f, T_g) \leq \frac{1}{2} \sup_{x, y \in X} |d_f(x, y) - d_g(x, y)| . \quad (4.84)$$

Following the logic of the proof of lemma 3.13, the distance between  $\pi_f(x)$  and  $\pi_f(y)$  is of the form

$$d_f(\pi_f(x), \tau) + d_f(\tau, \pi_f(y)) = f(x) - f(\tau) + f(y) - f(\tau) \quad (4.85)$$

where  $\tau$  is the lowest point of the geodesic path in  $T_f$  between  $\pi_f(x)$  and  $\pi_f(y)$ . This geodesic path on  $T_f$  admits preimages by  $\pi_f$  which are paths connecting  $x$  to  $y$ . These paths achieve the following supremum

$$\sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f \circ \gamma = f(\tau) \leq f(x) \wedge f(y) \quad (4.86)$$

where  $a \wedge b := \min\{a, b\}$  since by construction  $\gamma$  must always stay above  $f(\tau)$  and since for  $r > f(\tau)$ ,  $x$  and  $y$  lie in different connected components of  $X_r$ . If  $\nu$  is the analogous vertex to  $\tau$  on  $T_g$ , we have that

$$\begin{aligned} d_{GH}(T_f, T_g) &\leq \frac{1}{2} \sup_{x, y \in X} |d_f(x, y) - d_g(x, y)| \\ &= \frac{1}{2} \sup_{x, y \in X} |f(x) - g(x) + f(y) - g(y) - 2f(\tau) + 2g(\nu)| \\ &\leq \|f - g\|_{L^\infty} + \sup_{x, y \in X} \left| \sup_{\gamma: x \rightarrow y} \inf_{t \in [0, 1]} f \circ \gamma - \sup_{\eta: x \rightarrow y} \inf_{t \in [0, 1]} g \circ \eta \right| \\ &\leq 2 \|f - g\|_{L^\infty} , \end{aligned} \quad (4.87)$$

as desired. ■

*Remark 4.7.* One can be tempted to establish a general inequality between  $d_{GH}$  and  $d_b$  since both of these distances are bounded by the  $L^\infty$ -norm. However, this is not possible.

Indeed, there is a simple counter-example to  $d_{GH} \geq d_b$ . To illustrate this consider two barcodes  $k[s, -\infty[$  and  $k[s + \varepsilon, -\infty[$ . The bottleneck distance between these two is clearly  $\geq \varepsilon$ . But supposing that the functions  $f$  and  $g$  generating these barcodes are such that  $f = g + \varepsilon$  the trees  $T_g$  and  $T_f$  are isometric, so  $d_{GH}(T_f, T_g) = 0 < \varepsilon \leq d_b(\mathcal{B}(f), \mathcal{B}(g))$ .

Conversely, there are also counter-examples to  $d_b \geq d_{GH}$ , as this inequality would imply that two trees which have the same barcode are isometric. This is clearly false, as one can “glue” the bars of a given barcode in many different ways to give a tree, which generically will not be isometric.

## 5 A solution to the inverse problem

An interesting question is whether every (compact) tree stems from a function  $f : X \rightarrow \mathbb{R}$ . We can positively answer this question under the assumptions that  $\overline{\dim} T < \infty$  and that  $X = [0, 1]$  by constructing a function  $f : [0, 1] \rightarrow \mathbb{R}$ . The rest of this section will focus on proving the following theorem:

**Theorem 5.1.** *Let  $T$  be a compact  $\mathbb{R}$ -tree such that  $\overline{\dim} T < \infty$ . Then, for any  $\delta > 0$  it is possible to construct a continuous function  $f : [0, 1] \rightarrow \mathbb{R}$  of finite  $(\overline{\dim} T + \delta)$ -variation such that  $T = T_f$ . In particular, up to a reparametrization,  $f$  can be taken to be  $\frac{1}{\overline{\dim} T + \delta}$ -Hölder continuous.*

The idea is to once again use  $\varepsilon$ -simplifications  $T^\varepsilon$  for which we can construct a function by taking the contour of the tree. Such a construction is referred to as the Dyck path in the terminology of [28].

In what will follow, we will lay down notation which will simplify our task. In so doing, we will have shown the result for finite  $\mathbb{R}$ -trees. We then show the result for infinite trees.

### 5.1 Some preliminaries

We can regard a rooted discrete tree as being an operator with  $N$  inputs, where  $N$  is the number of leaves of the tree. There is a natural operation on the space of discrete trees which composes these operations by:

$$\text{Diagram 1} \circ \left( \begin{array}{c} \text{Diagram 2} \\ \text{Diagram 3} \end{array} \right) := \text{Diagram 4}$$

These objects are called **operads** and originated in the study of iterated loop spaces [5, 6, 21]. Since then, these objects have been studied in different fields for a variety of purposes [17, 19]. We will not give the explicit definition of an operad here, as we don't really need it, but we introduce this notion of composition of trees for notational simplicity.

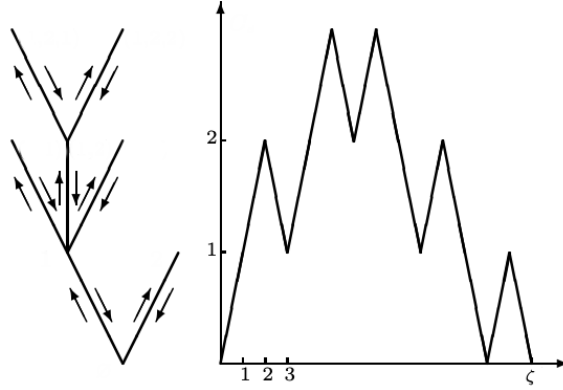


Figure 2: The Dyck path is the function  $f$  which assigns the height (the distance from the root) of each vertex of the tree as we wrap around the tree following a clockwise contour around it. There is a map  $\phi : T \mapsto [0, \zeta]$  where  $[0, \zeta]$  is now marked at the points at which  $f$  achieves its local maxima. The figure is taken from [15].

Given a discrete  $\mathbb{R}$ -tree  $T$ , if we have an embedding of  $T$  in  $\mathbb{R}^2$ , or equivalently, a partial order on its vertices, we can assign to  $T$  an interval  $I$  of a certain length with  $N$  marked points as well as a function  $f_T : I \rightarrow \mathbb{R}$ , where  $N$  is the number of leaves of  $T$ . Using the terminology of [28], a way to do this is by considering the so-called **Dyck path** or **contour path** where the path around  $T$  parametrized by arclength in  $T$ . The construction of the Dyck path has been carefully detailed in [15, 28], but it is better understood by looking at figure 2. By construction the equality:  $T_{f_T} = T$  holds for any discrete  $\mathbb{R}$ -tree  $T$ . Here, equality is taken up to isometry.

As per the description of figure 2, the construction of the Dyck path yields a map  $\phi$  which to  $T$  assigns an interval  $\phi(T)$  with  $N$  marked points. An example of the action of  $\phi$  is illustrated in figure 3.

This operation  $\phi$  is in fact a “morphism” with respect to a composition operation on the intervals, defined as follows. If we have an interval  $I$  with  $N$  marked points and  $N$  intervals  $J_k$  each with  $M_j$  marked points, the result of the operation  $I \circ (J_1, \dots, J_N)$  is the insertion of the marked interval  $J_k$  at the  $k$ th marked point of  $I$ . The length of  $I \circ (J_1, \dots, J_N)$  is

$$|I \circ (J_1, \dots, J_N)| = |I| + \sum_{k=1}^n |J_k|, \quad (5.88)$$



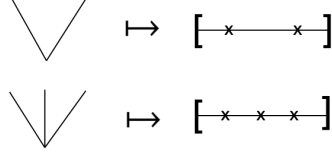


Figure 3: The action of  $\phi$  on trees with two and three leaves respectively. The length of the intervals assigned is exactly the length of the contour around the trees and the marked points are the points at which  $f_T$  achieves its maxima.

where  $|\cdot|$  denotes the lengths of the intervals. The fact that  $\phi$  is a “morphism” results from the definitions of compositions for trees and intervals. We can also define a variant of this morphism  $\phi$ , which we will call  $\phi_\lambda$ , which for any tree  $T$  simply scales the (marked) interval  $\phi(T)$  by a factor  $\lambda$ .

Given a tree  $T$  the Dyck path  $f_T : \phi(T) \rightarrow \mathbb{R}$  can be transformed into a function  $f_T^\lambda : \phi_\lambda(T) \rightarrow \mathbb{R}$  by setting

$$f_T^\lambda(x) := f_T(x/\lambda). \quad (5.89)$$

This is a rescaling of the  $x$ -axis which means that  $T_{f_T^\lambda} = T_{f_T} = T$  still holds. These equalities are taken up to isometry.

*Remark 5.2.* The definition of  $f_T^\lambda$  is readily generalizable to forests. If  $\mathcal{F}$  denotes a forest, then we define  $f_{\mathcal{F}}^\lambda = \bigsqcup_{T \in \mathcal{F}} f_T^\lambda$ .

For discrete trees, there is an upper bound of the number of vertices of the tree given its number of leaves.

**Lemma 5.3.** Let  $T$  be a rooted discrete tree,  $N$  be its number of leaves and  $V$  be its number of vertices, then

$$V \leq 2N - 1. \quad (5.90)$$

In particular, if the edges of  $T$  all have length 1, the contour of the tree can be done over an interval of length at most  $4N - 2$

*Proof.* For binary trees, it is known that [15, 28]

$$V = 2N - 1. \quad (5.91)$$

Given a tree with  $N$  leaves, we can obtain a binary tree with  $N$  leaves by simply blowing up the vertices which are non-binary. The inequality of the lemma follows. On a binary tree, the Dyck path passes through almost every point in  $T$  twice, so the length of the interval is exactly  $4N - 2$ . Since binary trees are the extremal case, a bound for all trees with  $N$  leaves follows. ■

The results above show the result of theorem 5.1 for finite trees, since their upper-box dimension is equal to 1.

## 5.2 Infinite trees: a construction idea

The concatenation of trees can be defined for  $\mathbb{R}$ -trees too in the obvious way. Given an infinite number of compositions, we can define a limit tree by defining it to be the limit of the partial compositions in the Gromov-Hausdorff sense. Ideally, we would like to have an equality of the following type

$$T = T^a \circ \overline{(T \setminus T^a)}, \quad (5.92)$$

where  $T \setminus T^a$  now denotes the rooted forest corresponding to the set  $T \setminus T^a$ . This equality is desirable because by taking infinitely many compositions, we can eventually recover the original tree  $T$ , by composing successive  $\varepsilon_n$ -simplifications with each other. However, this equality does not hold since  $T^a$  might not have the right amount of leaves for this operation to be well-defined. Nonetheless, we can decide to count the vertices  $T^a \cap \overline{(T \setminus T^a)}$  as leaves, so that the equality above holds.

In particular, the equality above would imply the following “proposition”.

**Fictional proposition 5.4.** *For any sequence  $(\varepsilon_n)_{n \in \mathbb{N}^*}$  such that  $\varepsilon_n \rightarrow 0$  monotonously, we have that for any compact  $\mathbb{R}$ -tree*

$$T = T^{\varepsilon_1} \circ (T^{\varepsilon_2} \setminus T^{\varepsilon_1}) \circ (T^{\varepsilon_3} \setminus T^{\varepsilon_2}) \circ \dots \quad (5.93)$$

For an infinite compact tree with  $\dim T < \infty$ , the idea is to take some appropriate rapidly decreasing sequence  $(\varepsilon_n)_{n \in \mathbb{N}^*}$  such that the interval

$$I = \phi_{\varepsilon_1}(T^{\varepsilon_1}) \circ \phi_{\varepsilon_2}(T^{\varepsilon_2} \setminus T^{\varepsilon_1}) \circ \phi_{\varepsilon_3}(T^{\varepsilon_3} \setminus T^{\varepsilon_2}) \circ \dots \quad (5.94)$$

has finite length. On each  $\phi_{\varepsilon_k}(T^{\varepsilon_k} \setminus T^{\varepsilon_{k-1}})$  we can consider the Dyck path on the forest  $T^{\varepsilon_k} \setminus T^{\varepsilon_{k-1}}$ . Defining a correct superposition of these Dyck paths, we would be done (*cf.* figure 4)

## 5.3 Infinite trees: the rigorous construction

Now that we have laid out the idea of the proof, we need to provide a rigorous construction of the ideas above. For this, we need to introduce multiple definitions.

**Definition 5.5.** Let  $I \subset \mathbb{R}_+$  be a marked interval with  $n$  marked points, which we will denote  $(i_k)_{1 \leq k \leq n}$ . Furthermore, let  $(J_k)_{1 \leq k \leq n}$  be a set of  $n$  marked intervals of  $\mathbb{R}_+$ , each with  $j_k$  marked points. Define  $\sigma_I : I \rightarrow I \circ (J_1, \dots, J_n)$  by

$$\sigma_I(x; J_1, \dots, J_n) := \left[ x + \sum_{i=1}^n \left[ \arg \max_k \{i_k < x\} \right] |J_i| \right] \in I \circ (J_1, \dots, J_n). \quad (5.95)$$

*Remark 5.6.* Fixing  $J_1, \dots, J_n$ ,  $\sigma_I$  is a bijective map onto its image, meaning every point  $y \in \sigma_I(I; J_1, \dots, J_n)$  admits a preimage in  $I$ , which we will denote by  $\sigma_I^{-1}(y; J_1, \dots, J_n)$ .

**Definition 5.7.** Let  $f : I \rightarrow \mathbb{R}$  be a continuous function from an interval  $I$  with  $n$  marked points and let  $(J_1, \dots, J_n)$  be intervals with each with  $j_i$  marked points as before. Abusing the notation, we define another function  $\sigma(-; J_1, \dots, J_n)$  which assigns a function on  $I$  to a function on  $\sigma_I(I; J_1, \dots, J_n)$  via the following formula

$$\sigma_I(f; J_1, \dots, J_n)(x) := \begin{cases} f(\sigma_I^{-1}(x; J_1, \dots, J_n)) & x \in \sigma_I(I; J_1, \dots, J_n) \\ \text{Linearly extend elsewhere} & \end{cases} \quad (5.96)$$

*Remark 5.8.* By continuity of  $f : I \rightarrow \mathbb{R}$ , this linear extension on  $I \circ (J_1, \dots, J_n)$  is in fact constant everywhere outside  $\sigma_I(I; J_1, \dots, J_n)$  (this is the dotted region in figure 4). Note also that  $\sigma_I(f; J_1, \dots, J_n)$  is continuous.

**Definition 5.9.** Given a tree  $T_f$  associated to a continuous function  $f$ , we define:

- The **projection onto the tree** as the mapping

$$\pi : X \rightarrow T_f = X/\{d_f = 0\}; \quad (5.97)$$

$$x \mapsto [f(x)] \quad (5.98)$$

- Let  $\tau \in T_f$ , define the **left preimage of  $\tau$** ,  $\overleftarrow{\tau}$  and the **right preimage of  $\tau$  by  $\pi$** ,  $\overrightarrow{\tau}$  as

$$\overleftarrow{\tau} := \inf \pi^{-1}(\tau) \quad (5.99)$$

$$\overrightarrow{\tau} := \sup \pi^{-1}(\tau). \quad (5.100)$$

**Definition 5.10.** Let  $T$  be a discrete rooted tree and  $T' \subset T$  be a subtree sharing roots with  $T$  and suppose that we have chosen some embeddings of  $T$  and  $T'$  on the plane such that these embeddings are consistent. Suppose there is a function  $f : I \rightarrow \mathbb{R}$  on a certain interval  $I$  such that  $T_f = T'$ . Then, the marking of  $I$  induced by  $T$  is the marking induced by marking the preimage  $\pi_f^{-1}(T' \cap \overline{(T \setminus T')})$  chosen in the following way:

- If  $\tau \in T' \cap \overline{(T \setminus T')}$  admits a single preimage, choose this preimage;
- Else, if the connected component of  $\tau$  in  $\overline{T \setminus T'}$  is smaller (with respect to the partial order on the tree induced by the embedding of  $T$ ) than every vertex strictly greater than  $\tau \in T'$ , choose  $\overleftarrow{\tau}$ . Otherwise, choose  $\overrightarrow{\tau}$ .

We will denote this marking operation by  $\mu(I; T', T, f)$ .

We can also define analogous maps to  $\sigma_I$ , but this time on the intervals  $J_k$  as follows.

**Definition 5.11.** Let  $I \subset \mathbb{R}_+$  be a marked interval with  $n$  marked points, which we will denote  $(i_k)_{\{1 \leq k \leq n\}}$ . Furthermore, let  $(J_k)_{\{1 \leq k \leq n\}}$  be a set of  $n$  marked intervals of  $\mathbb{R}_+$ , each with  $j_k$  marked points. Define  $\eta_I^{J_k} : J_k \rightarrow I \circ (J_1, \dots, J_n)$  by

$$\eta_I^{J_k}(x; J_1, \dots, J_n) := x + i_k. \quad (5.101)$$

These maps define a map  $\eta_I = \bigsqcup_k \eta_I^{J_k}$  on  $\bigsqcup_k J_k$  and  $\eta_I$  also induces a map on the functions  $f : \bigsqcup_k J_k \rightarrow \mathbb{R}$ , defined analogously to  $\sigma_I$ , which we shall also denote  $\eta_I$ .

With this notation, the construction is made in accordance to algorithm 2. A depiction of the mechanism of algorithm 2 can be found in figure 4.

---

**Algorithm 2:** Construction of approximants

---

**Output:** A set of unions of intervals  $(I_i)_{i \in \{1, \dots, n\}}$  and a set of functions on  $I_n$ ,  $(f_i : I_n \rightarrow \mathbb{R})_{i \in \{1, \dots, n\}}$

**Input:** An infinite tree  $T$  and  $a > 0$ .

$I_1 \leftarrow \phi(T^a)$  ;

$f_1 \leftarrow f_{T^a}$  ;

$I \leftarrow I_1$  ;

$i \leftarrow 1$  ;

**while**  $i \leq n$  **do**

$I_{i+1} := I_i \circ \phi_{\lambda^i}(\overline{T^{a/2^{i+1}} \setminus T^{a/2^i}})$  ;

$f \leftarrow \eta_{I_{i+1}}(f_{T^{a/2^{i+1}} \setminus T^{a/2^i}}; I_1, \dots, I_i)$  ;

$I_i \leftarrow \mu(I_i; T^{a/2^{i-1}}, T^{a/2^i}, f_i)$  ;

**for**  $j=1; j \leq i$  **do**

$I_j \leftarrow \sigma(I_j; \phi_{\lambda^i}(\overline{T^{a/2^{i+1}} \setminus T^{a/2^i}}))$  ;

$f_j \leftarrow \sigma(f_j; \phi_{\lambda^i}(\overline{T^{a/2^{i+1}} \setminus T^{a/2^i}}))$  ;

$j \leftarrow j + 1$  ;

**end**

$f_{i+1} := f_i + f$  ;

$i \leftarrow i + 1$  ;

**end**

**return**  $(I_i)_{i \in \{1, \dots, n\}}, (f_i)_{i \in \{1, \dots, n\}}$ .

---

For an infinite tree, it suffices to show that the sequence generated by this algorithm converges in the Gromov-Hausdorff sense to an interval of finite length  $I$  and that  $(f_i)_i$  converge in  $L^\infty(I)$  to some function  $f$ . The first thing we must show is thus:

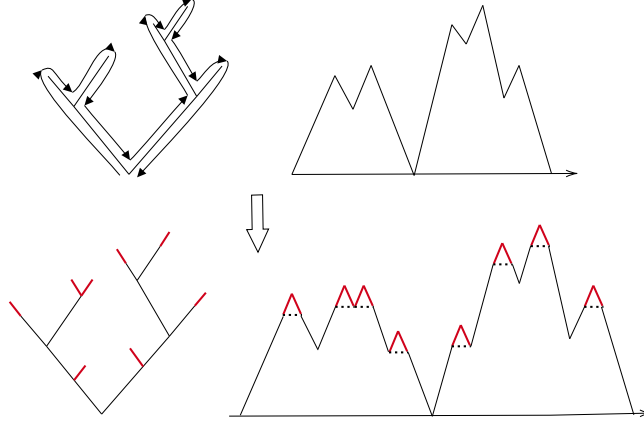


Figure 4: Starting from a tree  $T^{a/2^k}$  (black) we construct the Dyck path around it in the first step. Then, we look at  $T^{a/2^{k+1}}$  which leads to the addition of intervals (dotted), and a correction of the function at the  $k$ th step  $f_k$  (which is the function depicted in black, extended linearly over the new intervals). We can further define a function which by pasting the Dyck paths of the forest over the corresponding leaves, which leads to the function depicted in the second step (red and black).

**Lemma 5.12.** If  $T$  is a compact  $\mathbb{R}$ -tree of finite upper-box dimension, it is possible to define such an interval of finite length  $I$  defined by the construction above.

We need to show the convergence of the corresponding functions  $(f_n)_n$ . This can be done by proving that the sequence is Cauchy.

**Lemma 5.13.** Given the definition of functions  $f_n$  above, then the sequence  $(f_n)_{n \in \mathbb{N}^*}$  is Cauchy in  $L^\infty(I)$ , we have

$$\|f_n - f_m\|_{L^\infty} \leq a2^{-(n \wedge m)} \quad (5.102)$$

for any  $n$  and  $m \in \mathbb{N}^*$ .

By completeness of  $L^\infty(I)$ , the sequence  $(f_n)_{n \in \mathbb{N}^*}$  uniformly converges to a continuous function  $f$ . By virtue of stability theorem for trees (theorem 4.6) it follows that  $T$  is isometric to  $T_f$ . Using Picard's theorem (theorem 3.4)

$$\mathcal{V}(f) = \overline{\dim} T_f = \overline{\dim} T \quad (5.103)$$

which concludes the proof of theorem 5.1.

## 5.4 Proofs of the key lemmas

*Proof of lemma 5.12.* Recall that, according to the proof of theorem 3.7, the following equality holds for any tree  $T$

$$\overline{\lim}_{\varepsilon \rightarrow 0} \frac{\log N^\varepsilon}{\log(1/\varepsilon)} \vee 1 \leq \overline{\dim} T := \alpha. \quad (5.104)$$

Unpacking the definition of the lim sup, for any  $\delta > 0$  there is a  $a > 0$  such that for all  $\varepsilon < a$ , we have that

$$N^\varepsilon < \varepsilon^{-\alpha-\delta}. \quad (5.105)$$

Let us fix such a  $\delta$  and pick  $a$  small enough so that the condition above holds. For any  $n \in \mathbb{N}^*$ , the partial composition of intervals has length

$$|I_n| = |\phi(T^a)| + \sum_{k=1}^n \left| \phi_{\lambda^k}(T^{a/2^k} \setminus T^{a/2^{k-1}}) \right|. \quad (5.106)$$

However, we can bound  $\left| \phi_{\lambda^k}(T^{a/2^k} \setminus T^{a/2^{k-1}}) \right|$  by

$$\begin{aligned} \left| \phi_{\lambda^k}(T^{a/2^k} \setminus T^{a/2^{k-1}}) \right| &= \lambda^k \left| \phi(T^{a/2^k} \setminus T^{a/2^{k-1}}) \right| \\ &\leq \lambda^k \left( \frac{a}{2^k} \right) (4N^{a/2^k}), \end{aligned} \quad (5.107)$$

since on  $T^{a/2^k} \setminus T^{a/2^{k-1}}$  the distances between the vertices of each tree are at most  $a/2^k$  and there are at most  $4N^{a/2^k}$  such edges by virtue of lemma 5.3. Thus,

$$\left| \phi_{\lambda^k}(T^{a/2^k} \setminus T^{a/2^{k-1}}) \right| < 4\lambda^k \left( \frac{a}{2^k} \right)^{1-\alpha-\delta} = 4a^{1-\alpha-\delta} \left( 2^{\alpha+\delta-1}\lambda \right)^k. \quad (5.108)$$

Setting  $\lambda < 2^{1-\alpha-\delta}$   $I_n$  converges to some interval of finite length  $I$ , since the partial sums  $|I_n|$  converge.  $\blacksquare$

*Proof of lemma 5.13.* Suppose that  $n < m$ . It is sufficient to show that on  $I_m$  the equality holds, since in all further iterations of the algorithm, the functions  $f_n$  and  $f_m$  are locally constant over the intervals introduced. By definition of  $f_n$ ,  $f_n$  and  $f_m$  agree on  $I_n$ . Outside of this set,  $f_n$  is constant and the difference in the  $L^\infty$ -norm depends only on what happens above  $T^{a/2^n}$ , thus we can write

$$\|f_n - f_m\|_{L^\infty} \leq \left\| f_{T^{a/2^m} \setminus T^{a/2^n}} \right\|_{L^\infty} \quad (5.109)$$

by definition of  $f_n$ . However, the Dyck path on  $T^{a/2^m} \setminus T^{a/2^n}$  can at most reach a height of  $a(2^{-n} - 2^{-m}) < a2^{-n}$ , which finishes the proof.  $\blacksquare$

## 6 Limitations and prospects

One of the clear limitations of the work we have done so far is that, in the way it has been presented, trees only give a valid description of the  $H_0$ -barcode. It would be interesting to extend these notions and build an  $H_k$ -distance in general, which we expect would generate a forest instead of a tree. Furthermore, the metric invariants related to persistence diagrams in this work are only really useful in a  $C^0$ -setting. Indeed, if we have a more regular function  $f$  (say  $C^1$ ), then the tree  $T_f$  is finite, so the notions of dimension we have exploited in this paper carry no supplementary information (and so is the extension of Picard's theorem). In the spirit of Polterovitch *et al.*'s work, another interesting line of work would be to explore what kind of invariants are able to capture more on the regularity of Lipschitz functions and in particular, whether trees are useful in this setting or not. By virtue of connecting persistence theory and the study of trees, we have just gained access to a wide range of results readily available for stochastic processes.

## 7 Acknowledgements

The author would like to thank Pierre Pansu and Claude Viterbo for helping with the redaction of the manuscript as well as their guidance. Many thanks are also owed to Jean-François Le Gall and Nicolas Curien for the fruitful discussions without which some of this work would not have been possible.

## References

- [1] H. Adams, M. Aminian, E. Farnell, M. Kirby, J. Mirth, R. Neville, C. Peterson, and C. Shonkwiler. A fractal dimension for measures via persistent homology. *Abel Symposia*, pages 1–31, 2020.
- [2] H. Adams, T. Emerson, M. Kirby, R. Neville, C. Peterson, P. Shipman, S. Chepushtanova, E. Hanson, F. Motta, and L. Ziegelmeier. Persistence images: A stable vector representation of persistent homology. *Journal of Machine Learning Research*, 18(8):1–35, 2017.
- [3] R. J. Adler and J. E. Taylor. *Random Fields and Geometry*. Springer New York, 2007.
- [4] J.-M. Azaïs and M. Wschebor. *Level Sets and Extrema of Random Processes and Fields*. John Wiley & Sons, Inc., Jul 2008.
- [5] J. M. Boardman and R. M. Vogt. Homotopy-everything spaces. *Bulletin of the American Mathematical Society*, 74(6):1117–1123, nov 1968.

- [6] J. M. Boardman and R. M. Vogt. *Homotopy Invariant Algebraic Structures on Topological Spaces*. Springer Berlin Heidelberg, 1973.
- [7] D. Burago, Y. Burago, and S. Ivanov. *A course in metric geometry*, volume 33 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2001.
- [8] M. Carriere, S. Oudot, and M. Ovsjanikov. Sliced Wasserstein Kernel for Persistence Diagrams. In *ICML 2017 - Thirty-fourth International Conference on Machine Learning*, pages 1–10, Sydney, Australia, Aug. 2017.
- [9] F. Chazal, V. de Silva, M. Glisse, and S. Oudot. *The Structure and Stability of Persistence Modules*. Springer International Publishing, 2016.
- [10] F. Chazal and V. Divol. The density of expected persistence diagrams and its kernel based estimation. In B. Speckmann and C. D. Tóth, editors, *34th International Symposium on Computational Geometry (SoCG 2018)*, volume 99 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 26:1–26:15, Dagstuhl, Germany, 2018. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [11] I. Chiswell. *Introduction to  $\Lambda$ -Trees*. WORLD SCIENTIFIC, feb 2001.
- [12] D. Cohen-Steiner, H. Edelsbrunner, J. Harer, and Y. Mileyko. Lipschitz functions have  $L^p$ -stable persistence. *Foundations of Computational Mathematics*, 10(2):127–139, Jan 2010.
- [13] N. Curien, J.-F. Le Gall, and G. Miermont. The Brownian cactus i. scaling limits of discrete cactuses. *Ann. Inst. H. Poincaré Probab. Statist.*, 49(2):340–373, 05 2013.
- [14] V. Divol and T. Lacombe. Understanding the topology and the geometry of the persistence diagram space via optimal partial transport. *CoRR*, abs/1901.03048, 2019.
- [15] T. Duquesne and J.-F. Le Gall. *Random trees, Lévy processes and spatial branching processes*. Number 281 in *Astérisque*. Société mathématique de France, 2002.
- [16] T. Duquesne and J.-F. Le Gall. Probabilistic and fractal aspects of Lévy trees. *Probability Theory and Related Fields*, 131(4):553–603, Nov 2004.
- [17] V. Ginzburg and M. Kapranov. Koszul duality for operads. *Duke Mathematical Journal*, 76(1):203–272, oct 1994.
- [18] J.-P. Kahane. *Some random series of functions*. Cambridge University Press, Feb 1986.



- [19] J.-L. Loday. La renaissance des opérades. In *Séminaire Bourbaki. Volume 1994/95. Exposés 790-804*, pages 47–74, ex. Paris: Société Mathématique de France, 1996.
- [20] R. MacPherson and B. Schweinhart. Measuring shape with topology. *Journal of Mathematical Physics*, 53(7):073516, Jul 2012.
- [21] J. P. May. *The Geometry of Iterated Loop Spaces*. Springer Berlin Heidelberg, 1972.
- [22] Y. Mileyko, S. Mukherjee, and J. Harer. Probability measures on the space of persistence diagrams. *Inverse Problems*, 27(12):124007, Nov 2011.
- [23] S. Y. Oudot. *Persistence Theory - From Quiver Representations to Data Analysis*, volume 209 of *Mathematical surveys and monographs*. American Mathematical Society, 2015.
- [24] J. Picard. A tree approach to  $p$ -variation and to integration. *The Annals of Probability*, 36(6):2235–2279, Nov 2008.
- [25] I. Polterovich, L. Polterovich, and V. Stojisavljević. Persistence barcodes and laplace eigenfunctions on surfaces. *Geometriae Dedicata*, 201(1):111–138, Aug 2018.
- [26] L. Polterovich, D. Rosen, K. Samvelyan, and J. Zhang. Topological Persistence in Geometry and Analysis. *arXiv e-prints*, page arXiv:1904.04044, Apr 2019.
- [27] B. Schweinhart. Persistent homology and the upper box dimension. *Discrete & Computational Geometry*, Nov 2019.
- [28] R. P. Stanley. *Enumerative Combinatorics*. Cambridge University Press, 2009.
- [29] K. Turner, Y. Mileyko, S. Mukherjee, and J. Harer. Fréchet means for distributions of persistence diagrams. *Discrete & Computational Geometry*, 52(1):44–70, jul 2014.